

CSE544

Data Management

Lecture 9

Query Optimization – Part 1

Announcements

- Project proposal is due tonight
 - Submit via gitlab
 - Also update the title on the [spreadsheet](#).
- HW2 is due on Wednesday
- Wednesday: Kyle's lecture about HW3
- Review 3 due next Wednesday.

Query Optimization

Three major components:

1. Search space today
2. Cardinality and cost estimation Monday
3. Plan enumeration algorithms Wednesday

Search Space

Refers to the set of all possible alternate plans that the optimizers may explore

- **Access Path Selection**: how to use indices
- **Rewrite Rules**: what identities the optimizer applies

Access Path

Access path: implements a selection $\sigma_P(R)$

(P is sometimes called search argument SARG)

- A file scan, or
- An index *plus* a matching selection condition

Supplier(sid, sname, scity, sstate)

Access Path Selection

```
SELECT * FROM Supplier  
WHERE sid > 300  $\wedge$  scity='Seattle'
```

Indices:

B+-tree on **sid**; clustered

B+-tree on **scity**; unclustered

Supplier(sid, sname, scity, sstate)

Access Path Selection

```
SELECT * FROM Supplier  
WHERE sid > 300 ∧ scity='Seattle'
```

Indices:

B+-tree on **sid**; clustered

B+-tree on **scity**; unclustered

Which access path should we use?

Supplier(sid, sname, scity, sstate)

Access Path Selection

```
SELECT * FROM Supplier  
WHERE sid > 300  $\wedge$  scity='Seattle'
```

Indices:

B+-tree on **sid**; clustered

B+-tree on **scity**; unclustered

Which access path should we use? **Depends on stats**

Supplier(sid, sname, scity, sstate)

Access Path Selection

```
SELECT * FROM Supplier  
WHERE sid > 300  $\wedge$  scity='Seattle'
```

B(Supplier) = 100

T(Supplier) = 1000

V(Supplier,scity) = 20

Max(Supplier, sid) = 1000

Min(Supplier,sid) = 1

Indices:

B+-tree on **sid**; clustered

B+-tree on **scity**; unclustered

Which access path should we use? **Depends on stats**

Supplier(sid, sname, scity, sstate)

Access Path Selection

```
SELECT * FROM Supplier  
WHERE sid > 300  $\wedge$  scity='Seattle'
```

Indices:

B+-tree on **sid**; clustered

B+-tree on **scity**; unclustered

B(Supplier)

T(Supplier)

V(Supplier,scity) = 20

Max(Supplier, sid) = 1000

Min(Supplier,sid) = 1

Number of
distinct values
of scity

1000

Which access path should we use? **Depends on stats**

Supplier(sid, sname, scity, sstate)

Access Path Selection

```
SELECT * FROM Supplier
WHERE sid > 300 ∧ scity='Seattle'
```

$B(\text{Supplier}) = 100$

$T(\text{Supplier}) = 1000$

$V(\text{Supplier}, \text{scity}) = 20$

$\text{Max}(\text{Supplier}, \text{sid}) = 1000$

$\text{Min}(\text{Supplier}, \text{sid}) = 1$

Indices:

B+-tree on **sid**; clustered

B+-tree on **scity**; unclustered

Which access path should we use? **Depends on stats**

Supplier(sid, sname, scity, sstate)

Access Path Selection

```
SELECT * FROM Supplier
WHERE sid > 300 ∧ scity='Seattle'
```

$B(\text{Supplier}) = 100$

$T(\text{Supplier}) = 1000$

$V(\text{Supplier}, \text{scity}) = 20$

$\text{Max}(\text{Supplier}, \text{sid}) = 1000$

$\text{Min}(\text{Supplier}, \text{sid}) = 1$

Indices:

B+-tree on **sid**; clustered

B+-tree on **scity**; unclustered

Which access path should we use? **Depends on stats**

1. Sequential scan: cost = 100

Supplier(sid, sname, scity, sstate)

Access Path Selection

```
SELECT * FROM Supplier  
WHERE sid > 300  $\wedge$  scity='Seattle'
```

$B(\text{Supplier}) = 100$

$T(\text{Supplier}) = 1000$

$V(\text{Supplier}, \text{scity}) = 20$

$\text{Max}(\text{Supplier}, \text{sid}) = 1000$

$\text{Min}(\text{Supplier}, \text{sid}) = 1$

Indices:

B+-tree on **sid**; clustered

B+-tree on **scity**; unclustered

Which access path should we use? **Depends on stats**

1. Sequential scan: cost = 100
2. Index scan on **sid**: cost = $7/10 * 100 = 70$

Supplier(sid, sname, scity, sstate)

Access Path Selection

```
SELECT * FROM Supplier  
WHERE sid > 300  $\wedge$  scity='Seattle'
```

$B(\text{Supplier}) = 100$

$T(\text{Supplier}) = 1000$

$V(\text{Supplier}, \text{scity}) = 20$

$\text{Max}(\text{Supplier}, \text{sid}) = 1000$

$\text{Min}(\text{Supplier}, \text{sid}) = 1$

Indices:

B+-tree on **sid**; clustered

B+-tree on **scity**; unclustered

Which access path should we use? **Depends on stats**

1. Sequential scan: cost = 100
2. Index scan on **sid**: cost = $7/10 * 100 = 70$
3. Index scan on **scity**: cost = $1000/20 = 50$

Supplier(sid, sname, scity, sstate)

Access Path Selection

```
SELECT * FROM Supplier
WHERE sid > 300 ^ scity='Seattle'
```

$B(\text{Supplier}) = 100$

$T(\text{Supplier}) = 1000$

$V(\text{Supplier}, \text{scity}) = 20$

$\text{Max}(\text{Supplier}, \text{sid}) = 1000$

$\text{Min}(\text{Supplier}, \text{sid}) = 1$

Indices:

B+-tree on **sid**; clustered

B+-tree on **scity**; unclustered

Which access path should we use? **Depends on stats**

1. Sequential scan: cost = 100
2. Index scan on **sid**: cost = $7/10 * 100 = 70$
3. Index scan on **scity**: cost = $1000/20 = 50$

Optimal

Rewrite Rules

Search space is defined by the set of rewrite rules that the optimizer implements

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Example Optimization

```
SELECT x.sid, y.pno, y.quantity
FROM.  Supplier x, Supply y
WHERE x.sid = y.sid
      and x.scity = 'Seattle'
```

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Example Optimization

$\Pi_{x.sid, y.pno, y.quantity}$

$\sigma_{x.scity = 'Seattle'}$

$\bowtie_{x.sid = y.sid}$

Supplier x

Supply y

```
SELECT x.sid, y.pno, y.quantity
FROM Supplier x, Supply y
WHERE x.sid = y.sid
      and x.scity = 'Seattle'
```

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Push Selections Down

$\Pi_{x.sid, y.pno, y.quantity}$

$\sigma_{x.scity = 'Seattle'}$

$\bowtie_{x.sid = y.sid}$

Supplier x

Supply y

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Push Selections Down

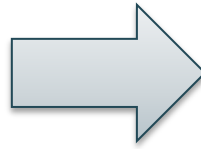
$\Pi_{x.sid, y.pno, y.quantity}$

$\sigma_{x.scity = 'Seattle'}$

$\bowtie_{x.sid = y.sid}$

Supplier x

Supply y



$\Pi_{x.sid, y.pno, y.quantity}$

$\sigma_{x.scity = 'Seattle'}$

$\bowtie_{x.sid = y.sid}$

Supplier x

Supply y

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Push Selections Down

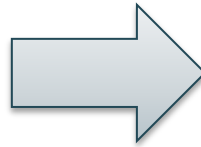
$\Pi_{x.sid, y.pno, y.quantity}$

$\sigma_{x.scity='Seattle'}$

$\bowtie_{x.sid = y.sid}$

Supplier x

Supply y



$\Pi_{x.sid, y.pno, y.quantity}$

$\sigma_{x.scity='Seattle'}$

$\bowtie_{x.sid = y.sid}$

Supplier x

Supply y

$$\sigma_C(R \bowtie S) = \sigma_C(R) \bowtie S \text{ when } C \text{ refers only to } R$$

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Push Selections Down

$\Pi_{x.sid, y.pno, y.quantity}$

$\sigma_{x.scity = 'Seattle' \text{ and } y.pno = 5}$

$\bowtie_{x.sid = y.sid}$

Supplier x

Supply y

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Push Selections Down

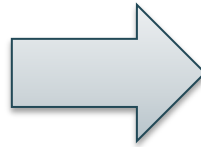
$\Pi_{x.sid, y.pno, y.quantity}$

$\sigma_{x.scity = 'Seattle' \text{ and } y.pno = 5}$

$\bowtie_{x.sid = y.sid}$

Supplier x

Supply y



$\Pi_{x.sid, y.pno, y.quantity}$

$\sigma_{x.scity = 'Seattle'}$

Supplier x

$\sigma_{y.pno = 5}$

Supply y

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Push Selections Down

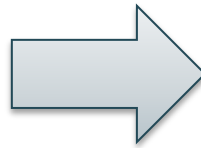
$\Pi_{x.sid, y.pno, y.quantity}$

$\sigma_{x.scity = 'Seattle' \text{ and } y.pno = 5}$

$\bowtie_{x.sid = y.sid}$

Supplier x

Supply y



$\Pi_{x.sid, y.pno, y.quantity}$

$\sigma_{x.scity = 'Seattle'}$

Supplier x

$\sigma_{y.pno = 5}$

Supply y

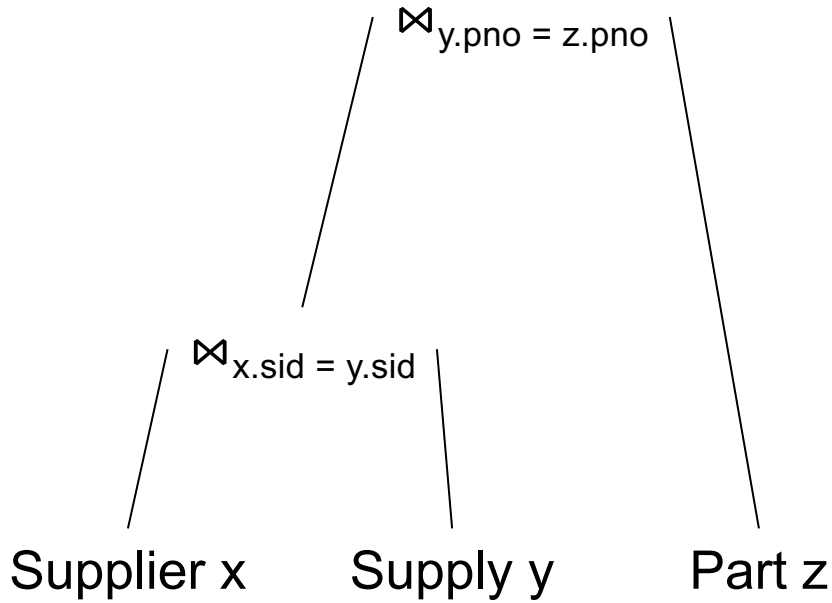
$$\sigma_{C1 \text{ and } C2}(R \bowtie S) = \sigma_{C1}(\sigma_{C2}(R \bowtie S)) = \sigma_{C1}(R \bowtie \sigma_{C2}(S)) = \sigma_{C1}(R) \bowtie \sigma_{C2}(S)$$

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Part(pno, pname, pprice)

Join Reorder

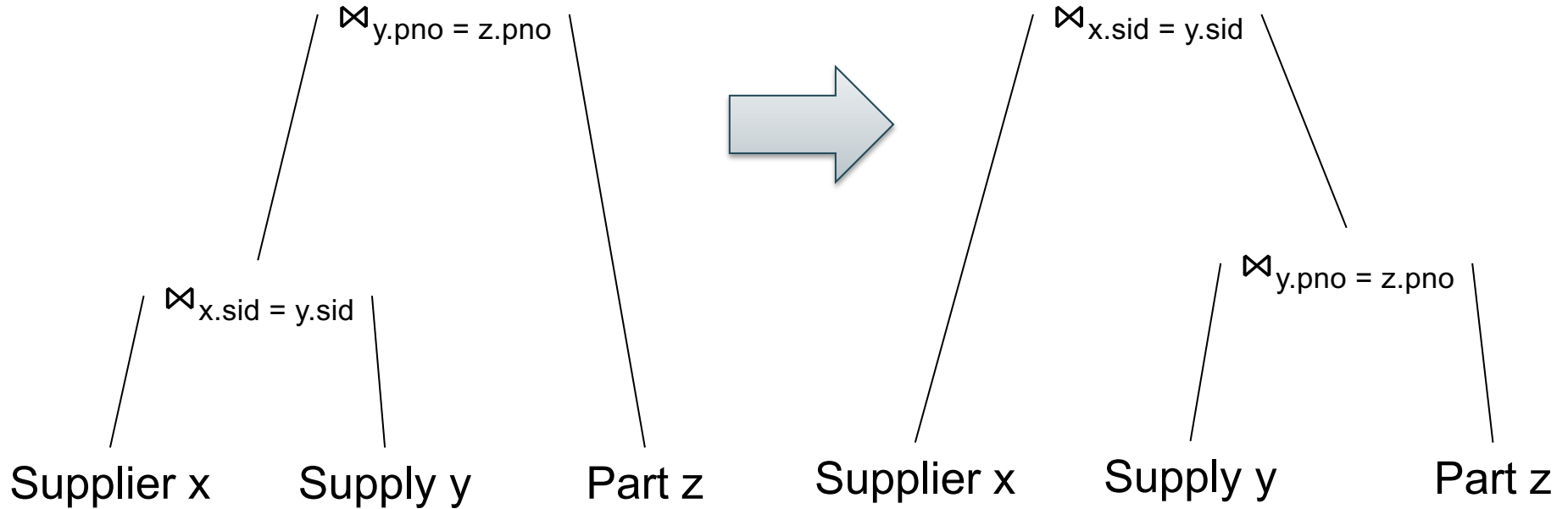


Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

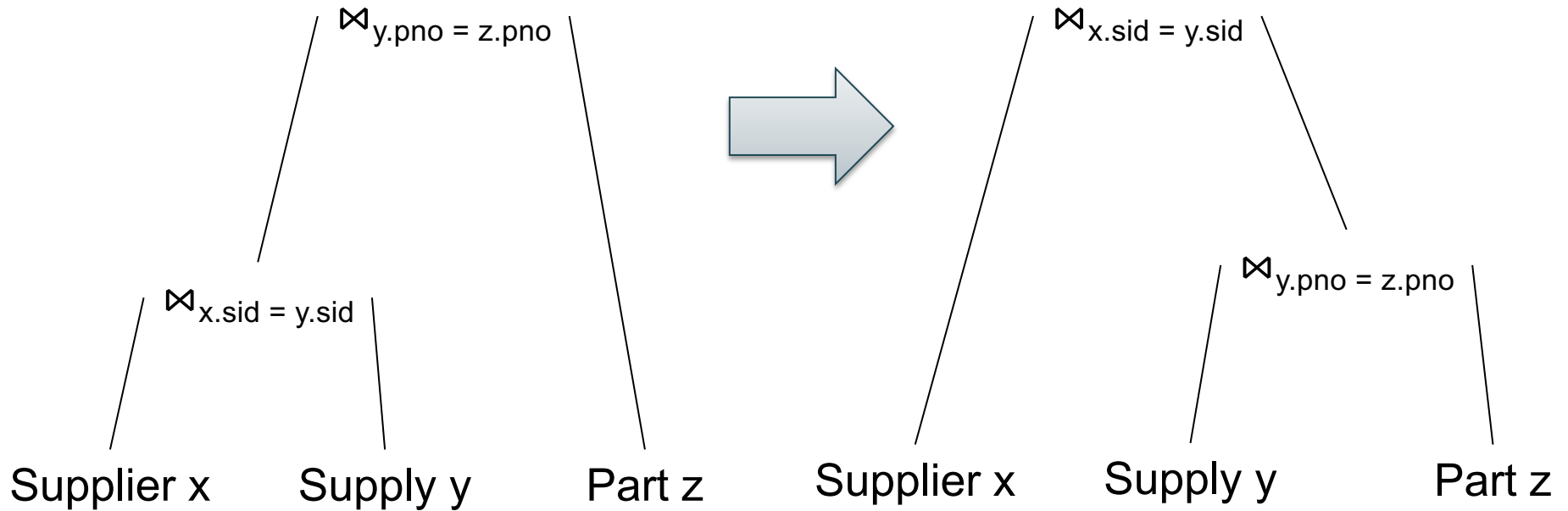
Part(pno, pname, pprice)

Join Reorder



Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)
Part(pno, pname, pprice)

Join Reorder



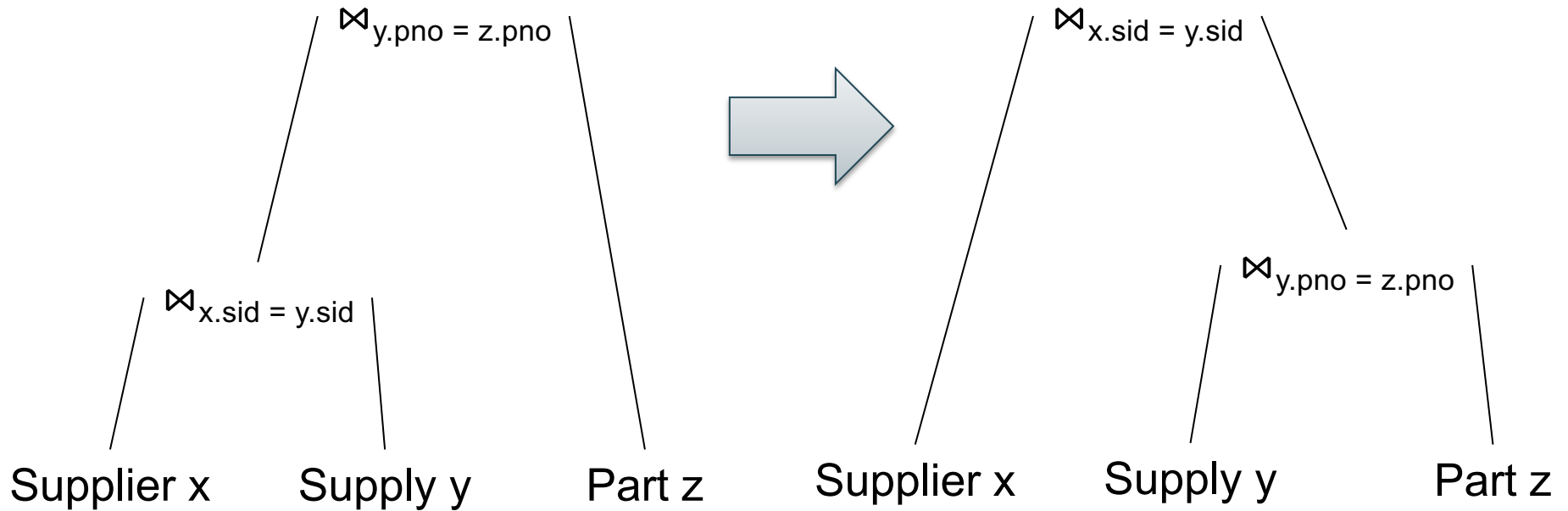
$$(R \bowtie S) \bowtie T = R \bowtie (S \bowtie T)$$

$$R \bowtie S = S \bowtie R$$

Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)
Part(pno, pname, pprice)

Join Reorder

Which plan is better?



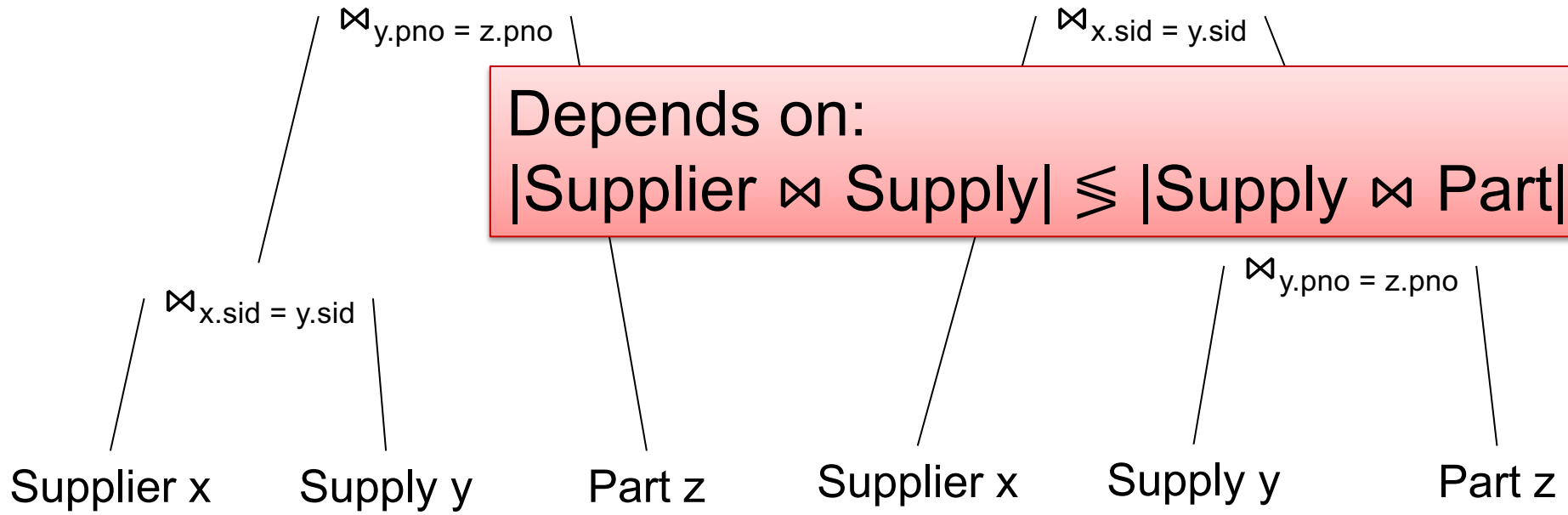
$$(R \bowtie S) \bowtie T = R \bowtie (S \bowtie T)$$

$$R \bowtie S = S \bowtie R$$

Supplier(sid, sname, scity, sstate)
 Supply(sid, pno, quantity)
 Part(pno, pname, pprice)

Join Reorder

Which plan is better?



Depends on:

$$|\text{Supplier} \bowtie \text{Supply}| \leq |\text{Supply} \bowtie \text{Part}|$$

$$(R \bowtie S) \bowtie T = R \bowtie (S \bowtie T)$$

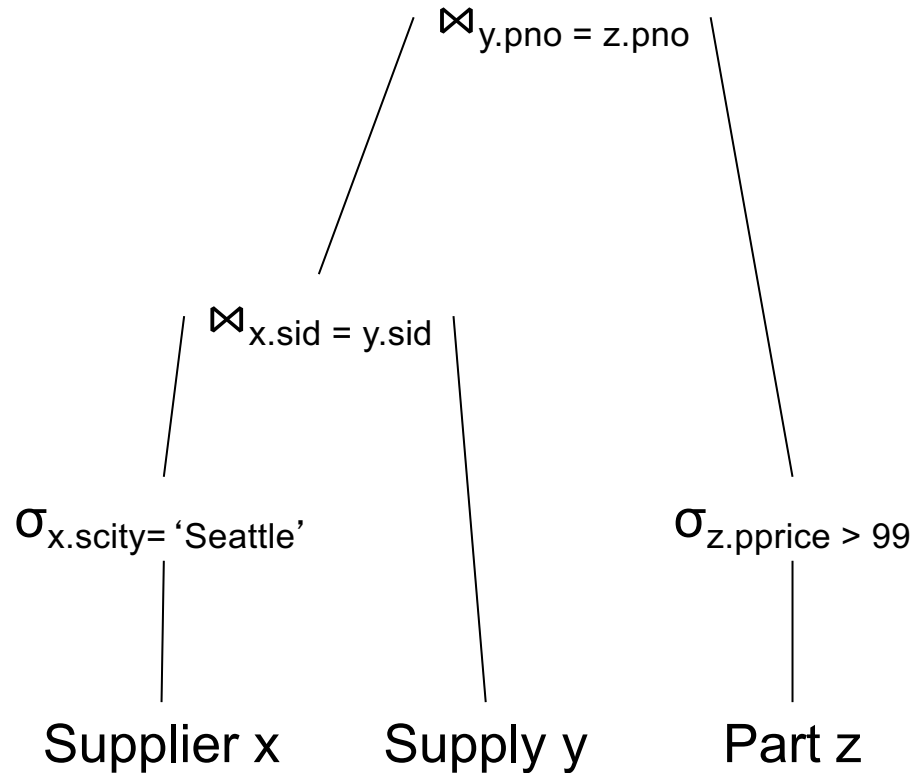
$$R \bowtie S = S \bowtie R$$

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Part(pno, pname, pprice)

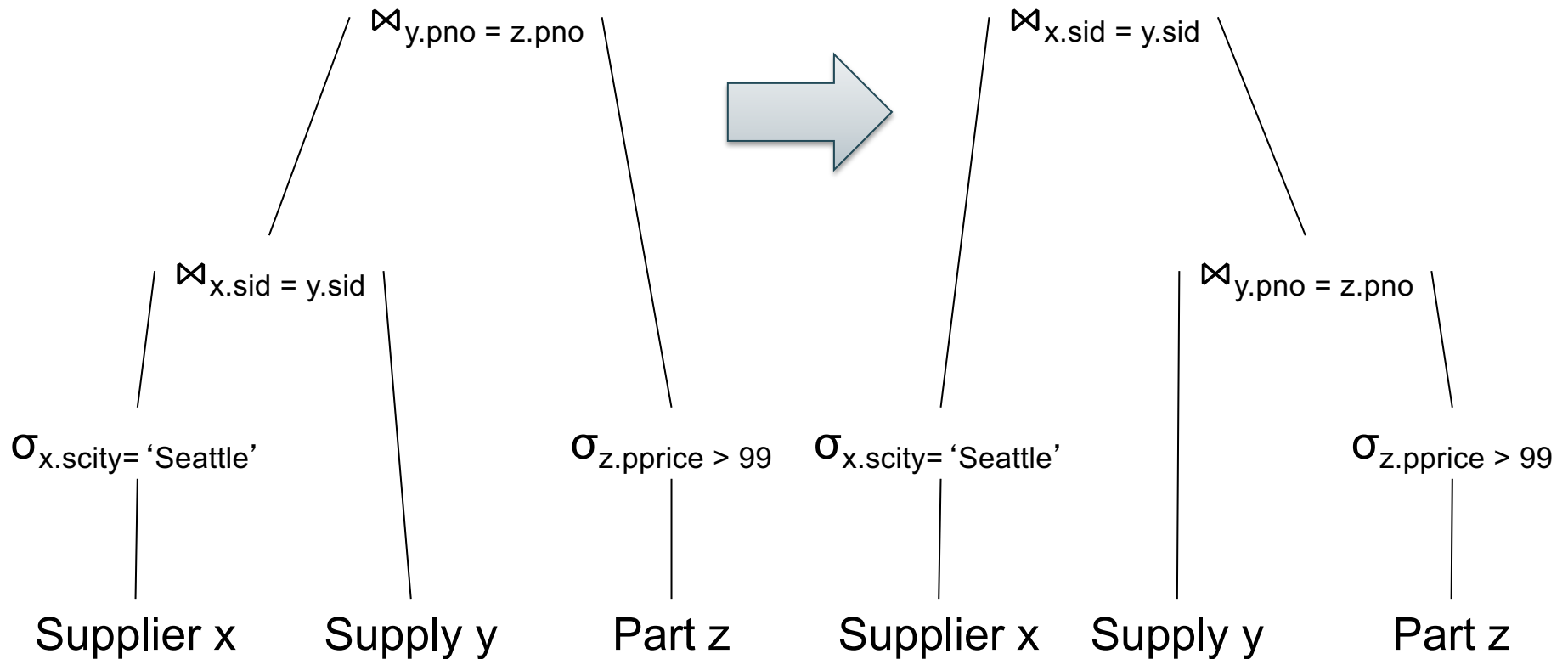
Join Reorder



Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)
Part(pno, pname, pprice)

Join Reorder

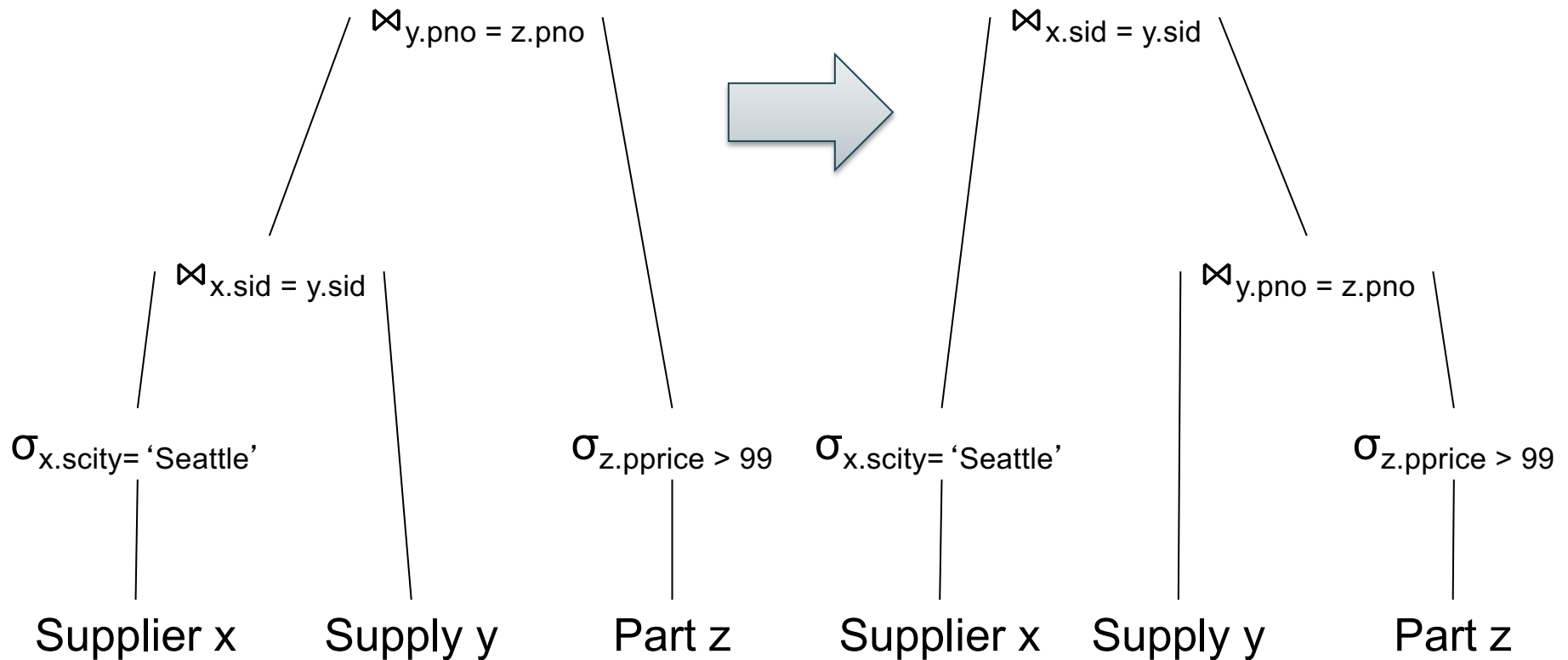
Which plan is better?



Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)
Part(pno, pname, pprice)

Join Reorder

Which plan is better?



Lesson: need sizes of $\sigma_{x.scity = 'Seattle'}$ (Supplier), $\sigma_{z.pprice > 99}$ (Part)

Discussion

- Both rewrite rules **and** cardinality estimations are needed for optimization
- They can be developed independently
- Next: a powerful optimization rule: aggregate pushdown or better known as generalized distributivity law

Motivation

- Try this in postgres

```
select count(*) from author;
```

Answer: 2652053

Time: 0.058 s

Motivation

- Try this in postgres

```
select count(*) from author;
```

Answer: 2652053

Time: 0.058 s

```
select count(*) from publication;
```

Answer: 5120896

Time: 0.062 s

Motivation

- Try this in postgres

```
select count(*) from author;
```

Answer: 2652053
Time: 0.058 s

```
select count(*) from publication;
```

Answer: 5120896
Time: 0.062 s

```
select count(*) from author, publication;
```

Timeout!!!

Motivation

- Try this in postgres

```
select count(*) from author;
```

Answer: 2652053
Time: 0.058 s

```
select count(*) from publication;
```

Answer: 5120896
Time: 0.062 s

```
select count(*) from author, publication;
```

Timeout!!!

But 3rd query is simply the **product** of the first two!

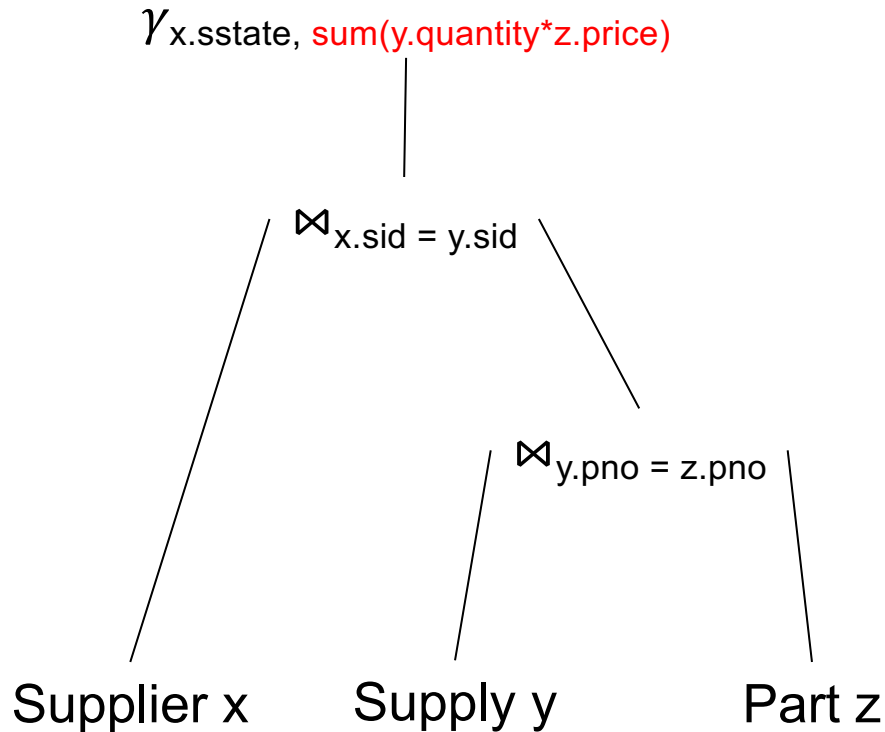
Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)
Part(pno, pname, pprice)

Aggregate Push-down

```
SELECT x.sstate, sum(y.quantity*z.price)
FROM Supplier x, Supply y, Part z
WHERE x.sid = y.sid and y.pno = z.pno
GROUP BY x.sstate
```

Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)
Part(pno, pname, pprice)

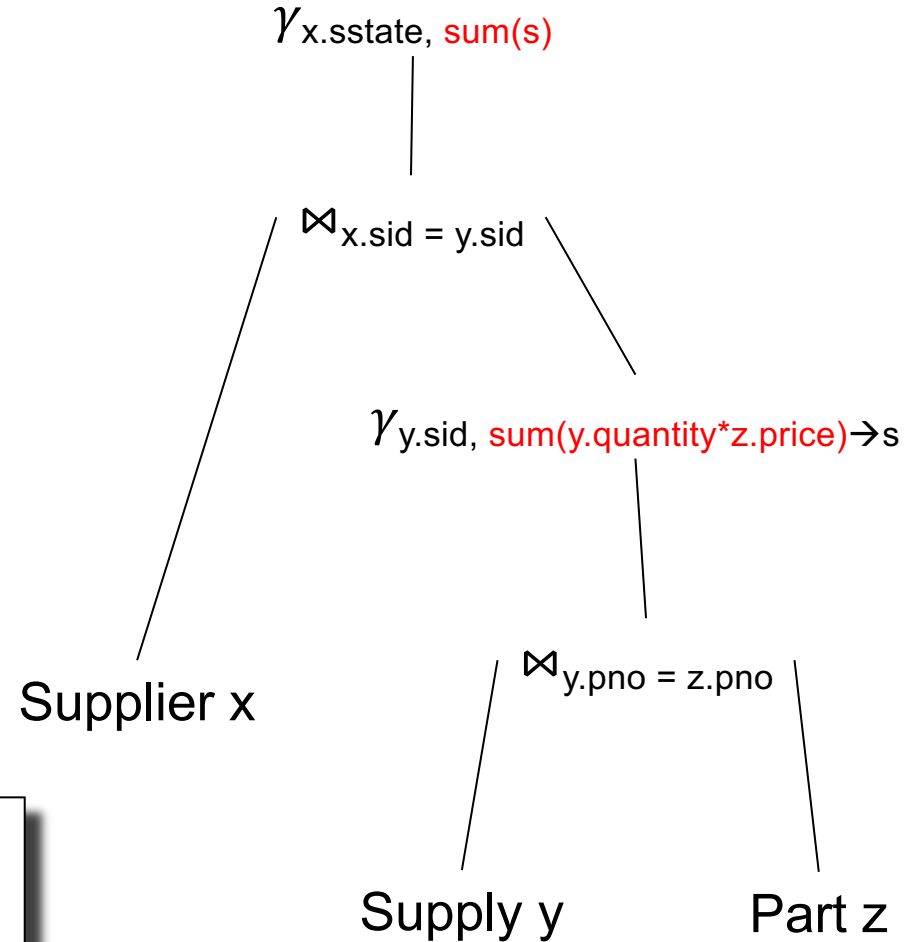
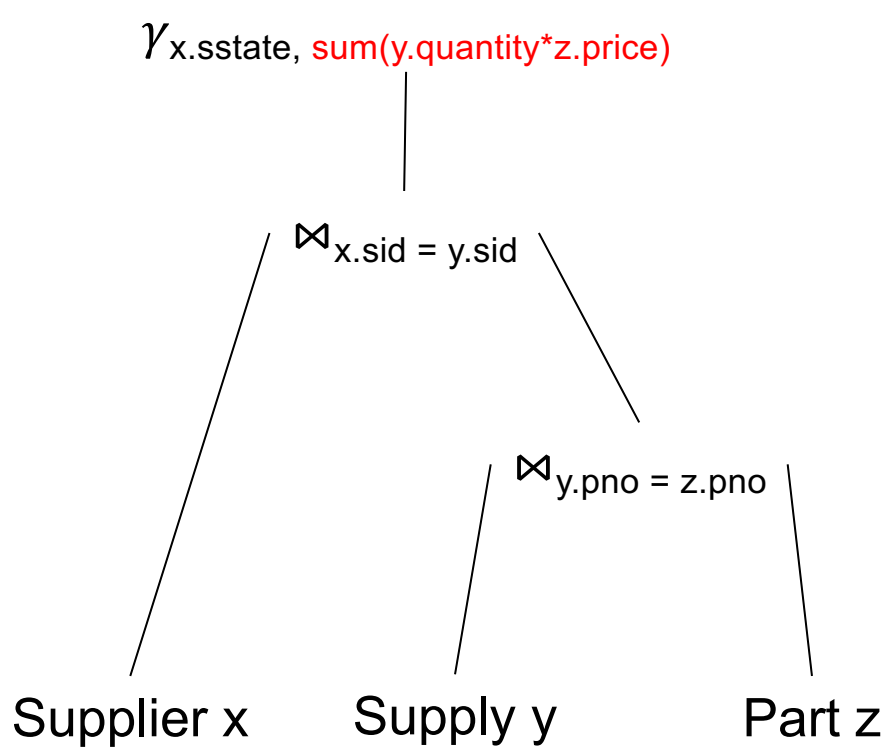
Aggregate Push-down



```
SELECT x.sstate, sum(y.quantity*z.price)
FROM Supplier x, Supply y, Part z
WHERE x.sid = y.sid and y.pno = z.pno
GROUP BY x.sstate
```

Supplier(sid, sname, scity, sstate)
 Supply(sid, pno, quantity)
 Part(pno, pname, pprice)

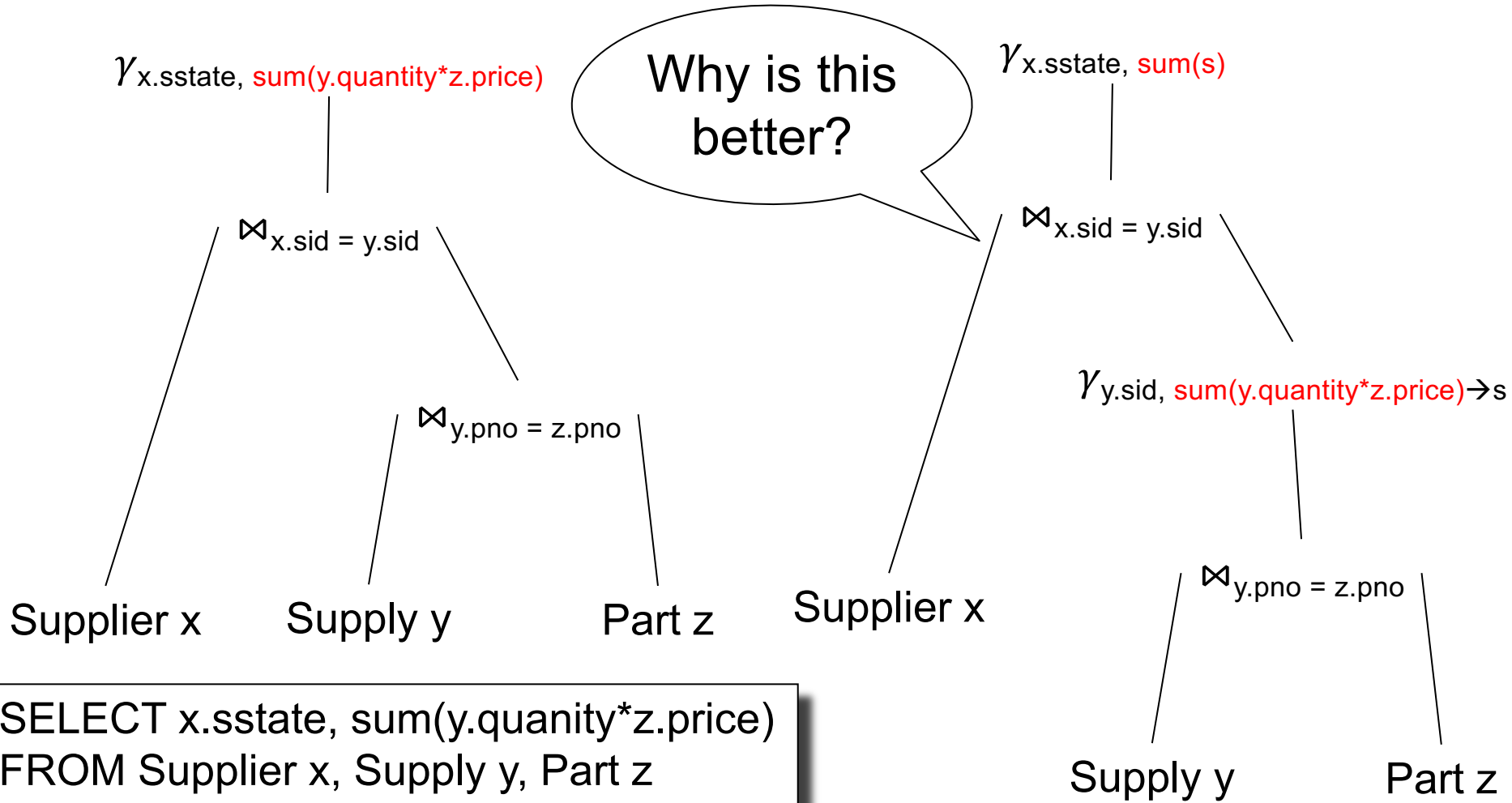
Aggregate Push-down



```
SELECT x.sstate, sum(y.quantity*z.price)
FROM Supplier x, Supply y, Part z
WHERE x.sid = y.sid and y.pno = z.pno
GROUP BY x.sstate
```


Supplier(sid, sname, scity, sstate)
 Supply(sid, pno, quantity)
 Part(pno, pname, pprice)

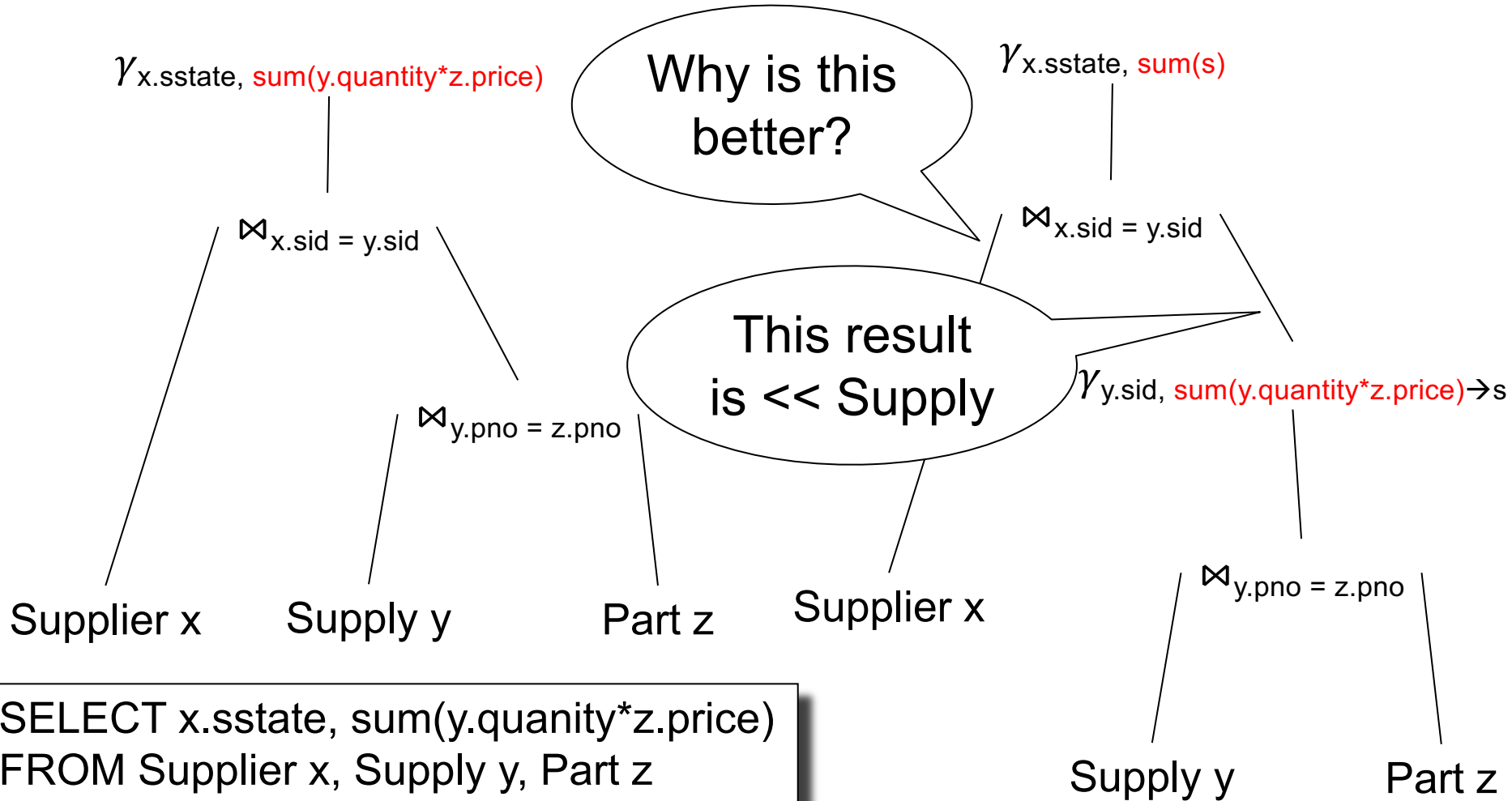
Aggregate Push-down



```
SELECT x.sstate, sum(y.quantity*z.price)
FROM Supplier x, Supply y, Part z
WHERE x.sid = y.sid and y.pno = z.pno
GROUP BY x.sstate
```

Supplier(sid, sname, scity, sstate)
 Supply(sid, pno, quantity)
 Part(pno, pname, pprice)

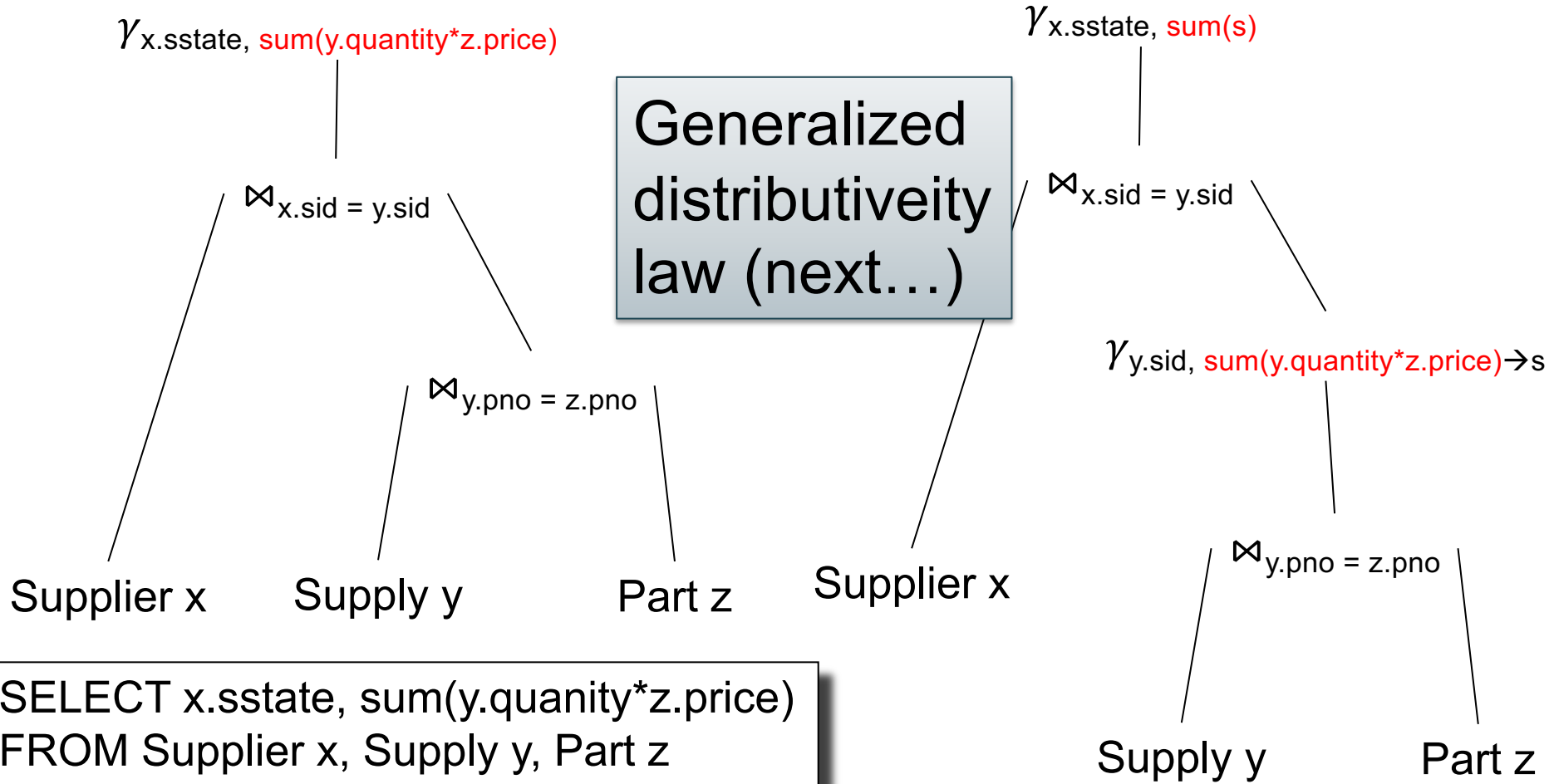
Aggregate Push-down



```
SELECT x.sstate, sum(y.quantity*z.price)
FROM Supplier x, Supply y, Part z
WHERE x.sid = y.sid and y.pno = z.pno
GROUP BY x.sstate
```

Supplier(sid, sname, scity, sstate)
 Supply(sid, pno, quantity)
 Part(pno, pname, pprice)

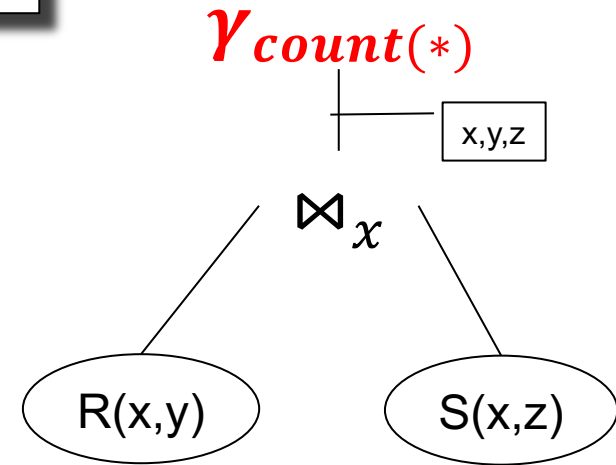
Aggregate Push-down



```
SELECT x.sstate, sum(y.quantity*z.price)
FROM Supplier x, Supply y, Part z
WHERE x.sid = y.sid and y.pno = z.pno
GROUP BY x.sstate
```

Generalized Distributivity Law

SELECT count(*) from R, S where R.x=S.x



Generalized Distributivity Law

SELECT count(*) from R, S where R.x=S.x

R:

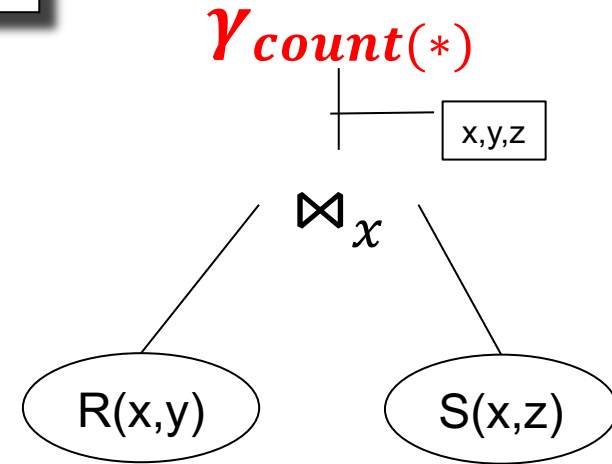
x	y
b	a
b	c
f	d
h	g

S:

x	z
b	g
b	k
h	m

Answer = 5

Runtime = ??????



Generalized Distributivity Law

SELECT count(*) from R, S where R.x=S.x

R:

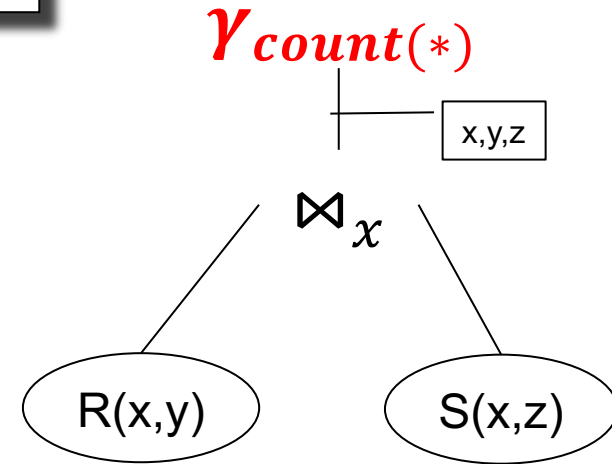
x	y
b	a
b	c
f	d
h	g

S:

x	z
b	g
b	k
h	m

Answer = 5

Runtime = $O(N^2)$



Generalized Distributivity Law

SELECT count(*) from R, S where R.x=S.x

R:

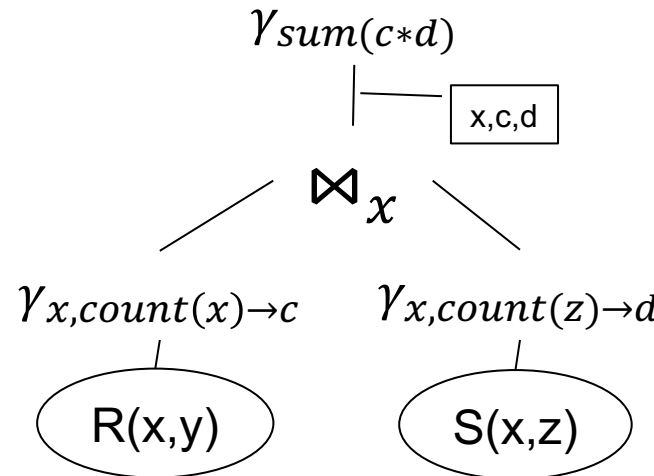
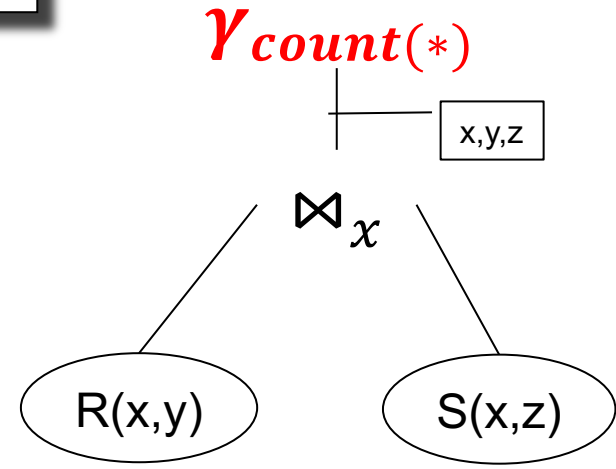
x	y
b	a
b	c
f	d
h	g

S:

x	z
b	g
b	k
h	m

Answer = 5

Runtime = $O(N^2)$



Generalized Distributivity Law

SELECT count(*) from R, S where R.x=S.x

R:

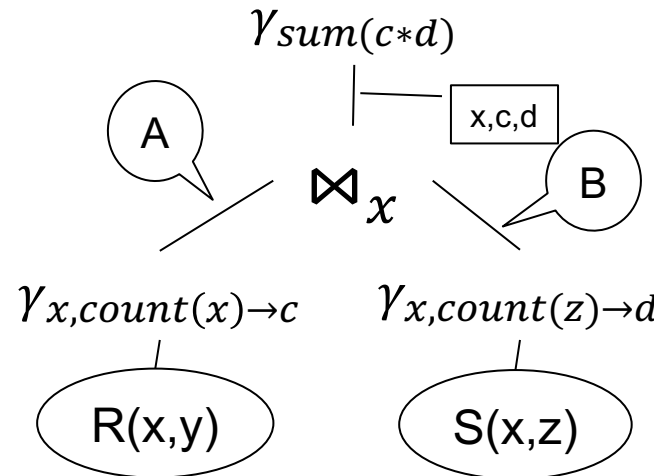
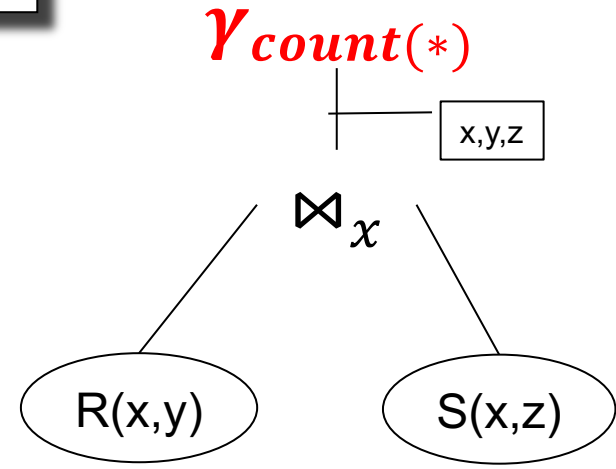
x	y
b	a
b	c
f	d
h	g

S:

x	z
b	g
b	k
h	m

Answer = 5

Runtime = $O(N^2)$



Generalized Distributivity Law

SELECT count(*) from R, S where R.x=S.x

R:

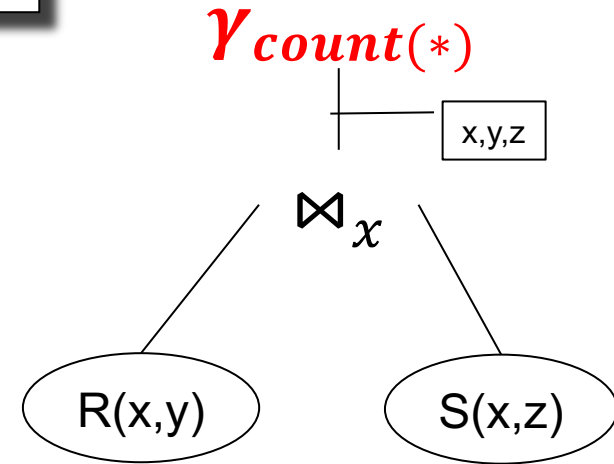
x	y
b	a
b	c
f	d
h	g

S:

x	z
b	g
b	k
h	m

Answer = 5

Runtime = $O(N^2)$



A:

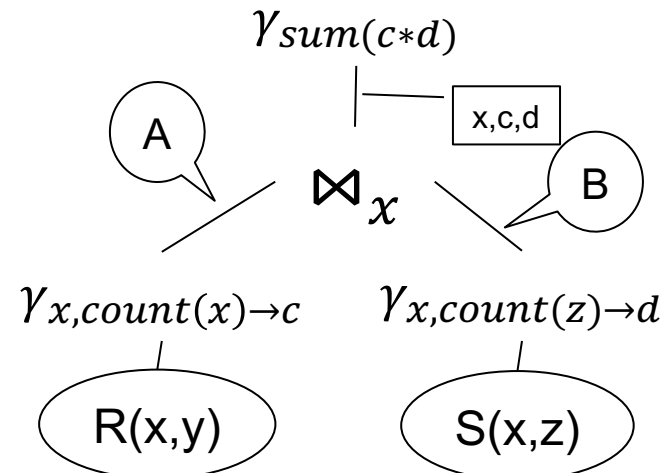
x	c
b	2
f	1
h	1

B:

x	d
b	2
h	1

$A \bowtie B$

x	c	d
b	2	2
h	1	1



Generalized Distributivity Law

SELECT count(*) from R, S where R.x=S.x

R:

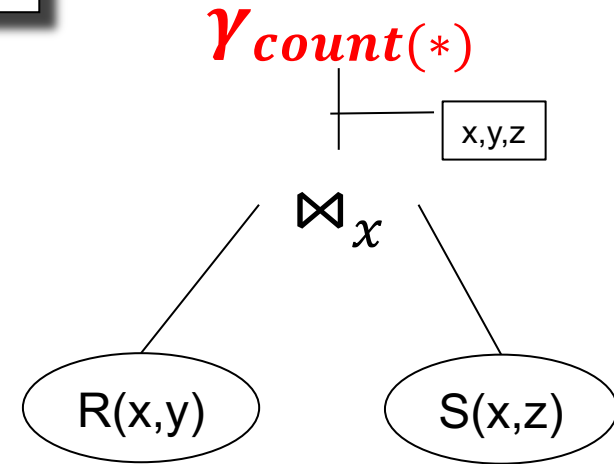
x	y
b	a
b	c
f	d
h	g

S:

x	z
b	g
b	k
h	m

Answer = 5

Runtime = $O(N^2)$



A:

x	c
b	2
f	1
h	1

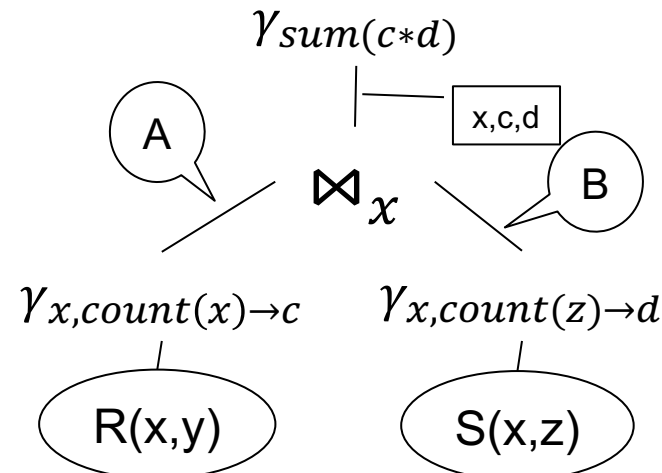
B:

x	d
b	2
h	1

$A \bowtie B$

x	c	d
b	2	2
h	1	1

Runtime = **?????**



Generalized Distributivity Law

SELECT count(*) from R, S where R.x=S.x

R:

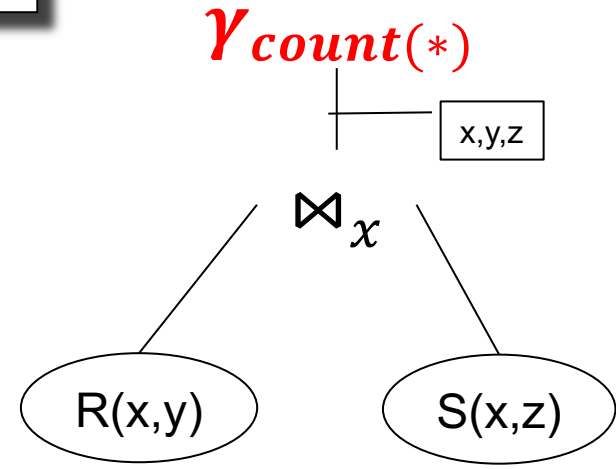
x	y
b	a
b	c
f	d
h	g

S:

x	z
b	g
b	k
h	m

Answer = 5

Runtime = $O(N^2)$



A:

x	c
b	2
f	1
h	1

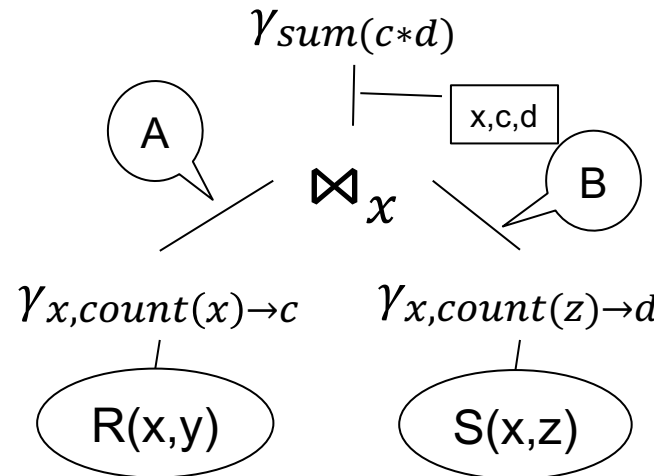
B:

x	d
b	2
h	1

$A \bowtie B$

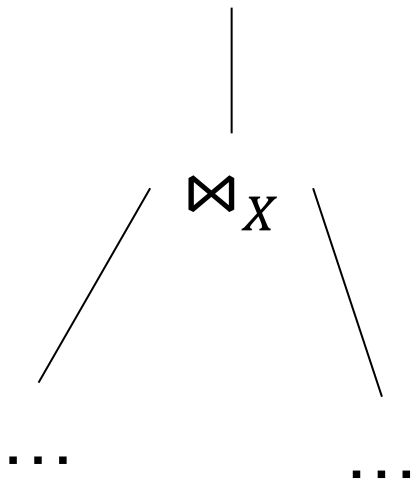
x	c	d
b	2	2
h	1	1

Runtime = $O(N)$



Generalized Distributivity Law

$\gamma_{Y,Z, \text{sum}}(A * B * C * \dots)$



Generalized Distributivity Law

$\gamma_{Y,Z, \text{sum}}(A * B * C * \dots)$

\bowtie_X

...

...

Y, A, C, E, ...
are here

Z, B, D, F, ...
are here

Generalized Distributivity Law

$\gamma_{Y,Z, \text{sum}}(A * B * C * \dots)$

\bowtie_X



\bowtie_X

?????????

...

...

...

...

Y, A, C, E, ...
are here

Z, B, D, F, ...
are here

Generalized Distributivity Law

$\gamma_{Y,Z,sum}(A*B*C*\dots)$

\bowtie_X

...

...

Y, A, C, E, ...
are here

Z, B, D, F, ...
are here



\bowtie_X

$\gamma_{X,Y,sum}(A*C*E\dots) \rightarrow S1$

...

...

Generalized Distributivity Law

$\gamma_{Y,Z,sum}(A*B*C*\dots)$

\bowtie_X

...

...

Y, A, C, E, ...
are here

Z, B, D, F, ...
are here



\bowtie_X

$\gamma_{X,Y,sum}(A*C*E\dots) \rightarrow S1$

...

??????????

...

Generalized Distributivity Law

$\gamma_{Y,Z,sum}(A*B*C*\dots)$

\bowtie_X

...

...

Y, A, C, E, ...
are here

Z, B, D, F, ...
are here



\bowtie_X

$\gamma_{X,Y,sum}(A*C*E\dots) \rightarrow S1$

$\gamma_{X,Z,sum}(B*D*F\dots) \rightarrow S2$

...

...

Generalized Distributivity Law

$\gamma_{Y,Z, \text{sum}}(A*B*C*\dots)$

\bowtie_X

...

...

Y, A, C, E, ...
are here

Z, B, D, F, ...
are here



???????

\bowtie_X

$\gamma_{X,Y, \text{sum}}(A*C*E\dots) \rightarrow S1$

$\gamma_{X,Z, \text{sum}}(B*D*F\dots) \rightarrow S2$

...

...

Generalized Distributivity Law

$\gamma_{Y,Z,sum}(A*B*C*\dots)$

\bowtie_X

...

...

Y, A, C, E, ...
are here

Z, B, D, F, ...
are here



$\gamma_{Y,Z,sum}(S1*S2)$

\bowtie_X

$\gamma_{X,Y,sum}(A*C*E\dots) \rightarrow S1$

$\gamma_{X,Z,sum}(B*D*F\dots) \rightarrow S2$

...

...

Discussion

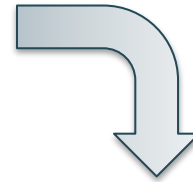
- Most DBMS support only some subset of the generalized distributivity law
- Fun activity: check if your favorite DBMS supports it
- Next: a rewrite rule that most DBMS support

Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Key / Foreign-Key

```
Select x.pno, x.quantity  
From Supply x, Supplier y  
Where x.sid = y.sid
```

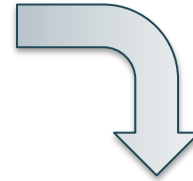


Supplier(sid, sname, scity, sstate)

Supply(sid, pno, quantity)

Key / Foreign-Key

```
Select x.pno, x.quantity  
From Supply x, Supplier y  
Where x.sid = y.sid
```

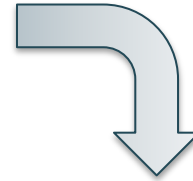


```
Select x.pno, x.quantity  
From Supply x
```

Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)

Key / Foreign-Key

```
Select x.pno, x.quantity  
From Supply x, Supplier y  
Where x.sid = y.sid
```



```
Select x.pno, x.quantity  
From Supply x
```

Only if these constraints hold:

1. Supplier.sid = key
2. Supply.sid = foreign key
3. Supply.sid NOT NULL

Summary: Rules

$$R \bowtie S = S \bowtie R$$

$$R \bowtie (S \bowtie T) = (R \bowtie S) \bowtie T$$

Summary: Rules

$$R \bowtie S = S \bowtie R$$

$$R \bowtie (S \bowtie T) = (R \bowtie S) \bowtie T$$

$$\sigma_{C_1 \wedge C_2}(R) = \sigma_{C_1}(\sigma_{C_2}(R))$$

Summary: Rules

$$R \bowtie S = S \bowtie R$$

$$R \bowtie (S \bowtie T) = (R \bowtie S) \bowtie T$$

$$\sigma_{C_1 \wedge C_2}(R) = \sigma_{C_1}(\sigma_{C_2}(R))$$

$$\sigma_C(R \bowtie S) = R \bowtie (\sigma_C(S)) \quad \text{if } C \text{ refers only to } S$$

Summary: Rules

$$R \bowtie S = S \bowtie R$$

$$R \bowtie (S \bowtie T) = (R \bowtie S) \bowtie T$$

$$\sigma_{C_1 \wedge C_2}(R) = \sigma_{C_1}(\sigma_{C_2}(R))$$

$$\sigma_C(R \bowtie S) = R \bowtie (\sigma_C(S)) \quad \text{if } C \text{ refers only to } S$$

$$\gamma_{X, \text{sum}(B)} \left(\gamma_{X, Y, \text{sum}(A) \rightarrow B}(R) \right) = \gamma_{X, \text{sum}(A)}(R)$$

Summary: Rules

$$R \bowtie S = S \bowtie R$$

$$R \bowtie (S \bowtie T) = (R \bowtie S) \bowtie T$$

$$\sigma_{C_1 \wedge C_2}(R) = \sigma_{C_1}(\sigma_{C_2}(R))$$

$$\sigma_C(R \bowtie S) = R \bowtie (\sigma_C(S)) \quad \text{if } C \text{ refers only to } S$$

$$\gamma_{X, \text{sum}(B)} \left(\gamma_{X, Y, \text{sum}(A) \rightarrow B}(R) \right) = \gamma_{X, \text{sum}(A)}(R)$$

In particular,
 $\delta(\delta(R)) = \delta(R)$

Summary: Rules

$$R \bowtie S = S \bowtie R$$

$$R \bowtie (S \bowtie T) = (R \bowtie S) \bowtie T$$

$$\sigma_{C_1 \wedge C_2}(R) = \sigma_{C_1}(\sigma_{C_2}(R))$$

$$\sigma_C(R \bowtie S) = R \bowtie (\sigma_C(S)) \quad \text{if } C \text{ refers only to } S$$

$$\gamma_{X, \text{sum}(B)} \left(\gamma_{X, Y, \text{sum}(A) \rightarrow B}(R) \right) = \gamma_{X, \text{sum}(A)}(R)$$

$$\gamma_{X, \text{sum}(A * B)}(R \bowtie_{Y=Z} S) = \\ \gamma_{X, \text{sum}(C * D)} \left(\gamma_{X, Y, \text{sum}(A) \rightarrow C}(R) \bowtie_{Y=Z} \gamma_{X, Z, \text{sum}(B) \rightarrow D}(S) \right)$$

In particular,
 $\delta(\delta(R)) = \delta(R)$

Summary: Rules

$$R \bowtie S = S \bowtie R$$

$$R \bowtie (S \bowtie T) = (R \bowtie S) \bowtie T$$

$$\sigma_{C_1 \wedge C_2}(R) = \sigma_{C_1}(\sigma_{C_2}(R))$$

$$\sigma_C(R \bowtie S) = R \bowtie (\sigma_C(S)) \quad \text{if } C \text{ refers only to } S$$

$$\gamma_{X, \text{sum}(B)} \left(\gamma_{X, Y, \text{sum}(A) \rightarrow B}(R) \right) = \gamma_{X, \text{sum}(A)}(R)$$

$$\gamma_{X, \text{sum}(A * B)}(R \bowtie_{Y=Z} S) = \gamma_{X, \text{sum}(C * D)} \left(\gamma_{X, Y, \text{sum}(A) \rightarrow C}(R) \bowtie_{Y=Z} \gamma_{X, Z, \text{sum}(B) \rightarrow D}(S) \right)$$

$$\Pi_{\text{attr}(R)}(R \bowtie_{X=Y} S) = R \quad \text{if } S.Y \text{ key, } R.X \text{ fk not null}$$

...

In particular,
 $\delta(\delta(R)) = \delta(R)$

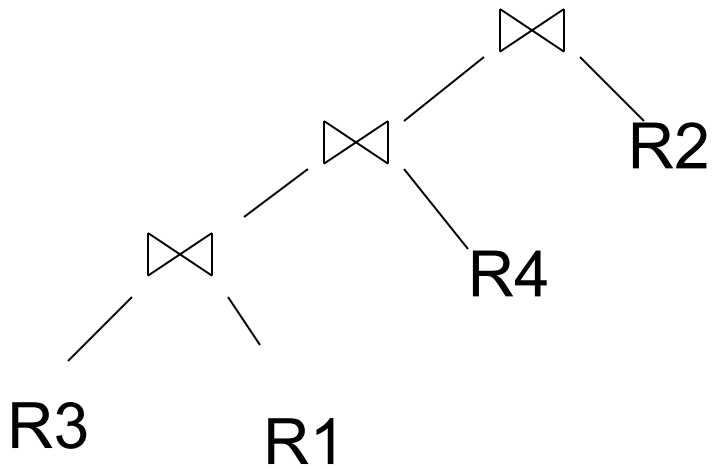
Practice

- Database optimizers typically have a database of rewrite rules
- E.g. SQL Server: about 500 rules
- Rules become complex as they need to serve specialized types of queries

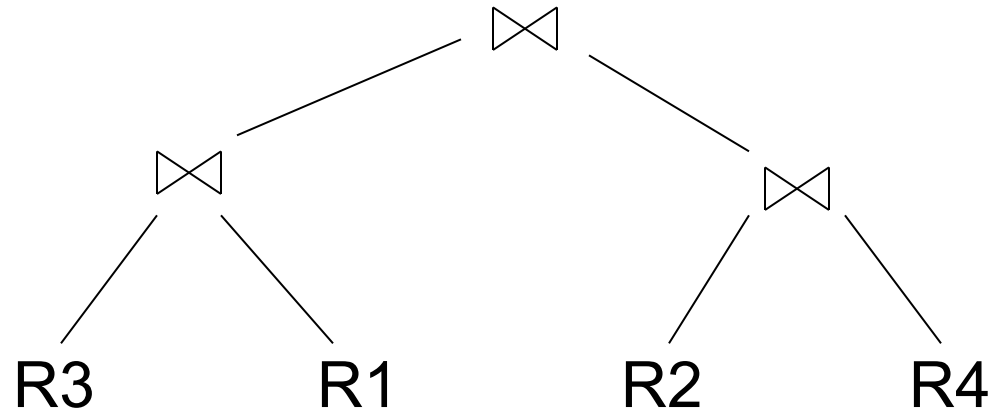
Search Space Challenges

- Search space is huge
- Typical compromises:
 - Left-deep plans
 - Plans without cartesian products

Left-Deep Plans and Bushy Plans



Left-deep plan



Bushy plan

Cartesian Product

- Plan without cartesian product:
 $(R(A, B) \bowtie S(B, C)) \bowtie T(C, D)$
- Plan with cartesian product
 $(R(A, B) \bowtie T(C, D)) \bowtie S(B, C)$
- Some optimizers avoid cartesian products

Query Optimization

Three major components:

1. Search space today
2. Cardinality and cost estimation Monday
3. Plan enumeration algorithms Wednesday